

HW2: Gradient Boosting

HW 2: Gradient Boosting

- Given a training data A, a testing data B, a target **testing** accuracy C
- Train a gradient boosting model and output predictions
- 10 test cases. Each case weights 1 pt.
- The compilation is considered failed if it does not finish in **5 minute**.
- A test case is considered **incorrect** if it does not finish in **2 minutes**.
- **Correct GPU solutions will get 5 pts bonus.**
- The **summation** of the execution time across 10 cases will be used to rank **correct** solutions.

Testing Environment

- `ssh yourusername@granger.cs.rit.edu`
- Intel(R) Xeon(R) CPU E5-2650 v4 @ 2.20GHz
- 48 threads in total (2 sockets, 12 cores per socket, 2 threads per core)
- 251 GB memory
- GPU: Tesla P4
- Testing limit:
 - 8 threads `taskset -c`
 - 1 GPU

Output Format

- N lines
- Each line contains an integer
 - The predicted class for each instance

What Do We Need to Do?

- We are required to complete two scripts
- `compiler.sh`
 - it is executed once before the actual testing starts
- `run.sh`
 - it should takes 4 arguments
 - the first argument is the training data file name
 - the second argument is the testing data file name
 - the third argument is the target testing accuracy
 - the fourth one is the file name that you should write your results into

HW2

- `bunzip2 mnist.bz2`
- `bunzip2 mnist.t.bz2`

- `bash run.sh <train> <test> <acc> <out>`
 - `bash run.sh mnist.t mnist.t 0.9 sample1.out`
 - `bash run.sh mnist mnist 0.9 sample2.out`
 - `bash run.sh mnist mnist.t 0.9 sample2.out`
 - ...
 - `bash run.sh mnist mnist.t 0.97 sample10.out`
- We guarantee that testing data will not have more features than the training data.

Neural Network Preliminary Tensor Computing

Weijie Zhao

02/17/2026

- Scalar
- Vector
- Matrix
- Tensor

Neural Network Preliminary Tensor Computing

- Rank
- Dimension

Weijie Zhao

02/17/2026

- Scalar
- Vector
- Matrix
- Tensor

Neural Network Preliminary

~~Tensor~~ Computing

Matrix

- Rank
- Dimension

Weijie Zhao

02/17/2026

- Matrix multiplication
- Non-linear activation
- Gradient descent

•Scalar

•Vector

•Matrix

•Tensor

•Rank

•Dimension

Neural Network Preliminary

~~Tensor~~ Computing Matrix

Weijie Zhao

02/17/2026

- Matrix multiplication
- Non-linear activation
- ~~Gradient descent~~ Graduate student descent

- Scalar
- Vector
- Matrix
- Tensor
 - Rank
 - Dimension

Neural Network Preliminary

~~Tensor~~ Computing

Matrix

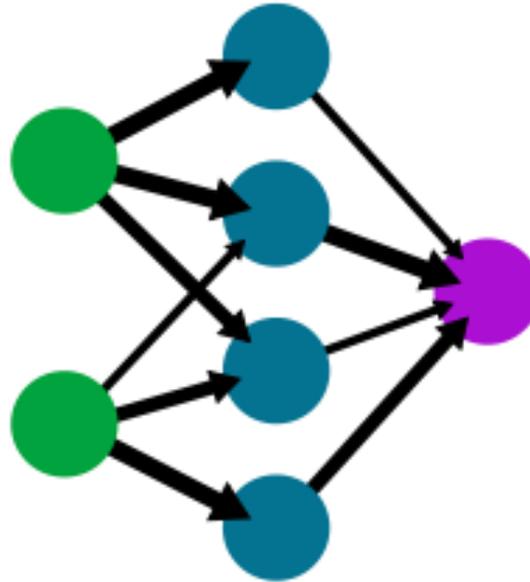
Weijie Zhao

02/17/2026

Neural Networks

A simple neural network

input layer hidden layer output layer



Deep Learning Framework Implementation

- Knowing the things under the hood
- Deployment
- Deployment on emerging hardware

Deep Learning Language

- How to represent a deep neural network?

Deep Learning Language

- How to represent a deep neural network?
 - Abstraction

Deep Learning Language

- How to represent a deep neural network?
 - Abstraction
- How to implement/deploy the abstraction deep neural network?

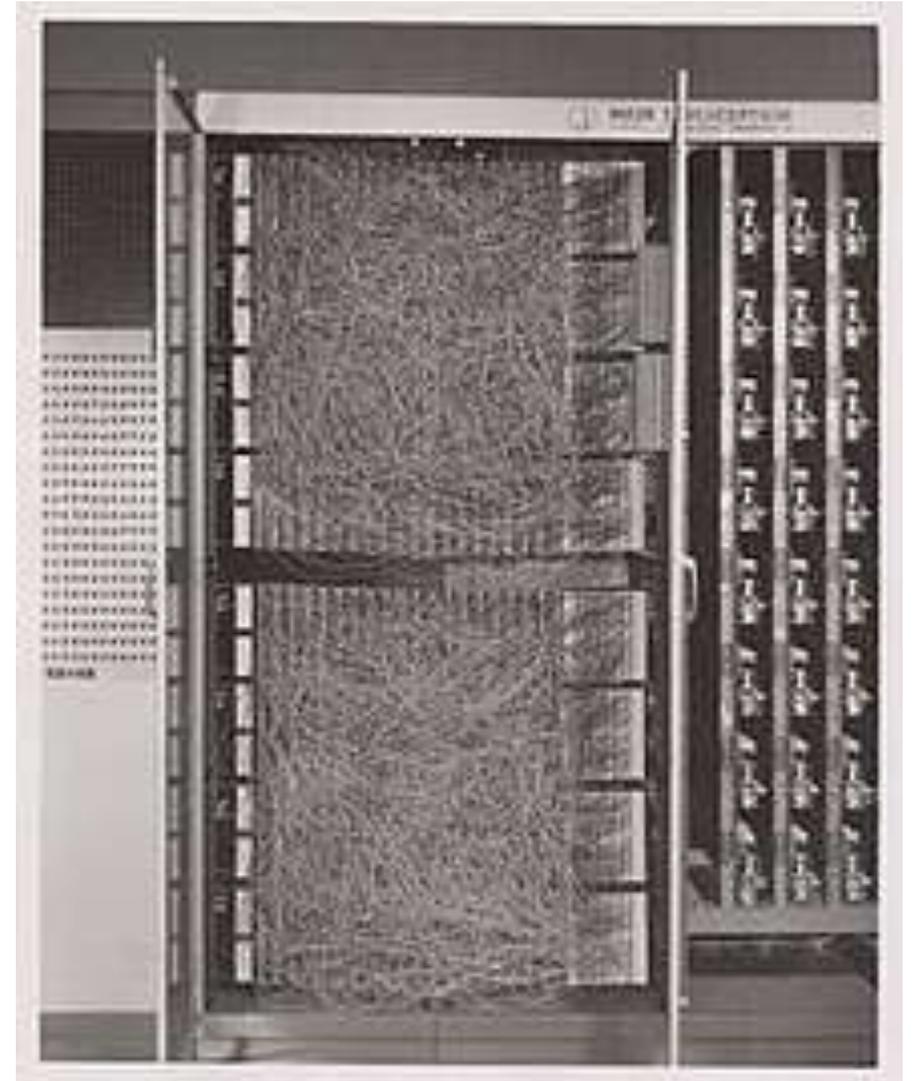
Deep Learning Language

- How to represent a deep neural network?
 - Abstraction
- How to implement/deploy the abstraction deep neural network?
 - Build tensor operations workflow
 - Implement high-performance low-level operations

Perceptron

- The perceptron was invented in 1943 by McCulloch and Pitts.
- The first implementation was a machine built in 1958 at the Cornell Aeronautical Laboratory by Frank Rosenblatt

$$f(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{w} \cdot \mathbf{x} + b > 0, \\ 0 & \text{otherwise} \end{cases}$$



Perceptron

- Linear Layer
 - Matrix multiplication
 - Addition
- ReLU

Tensor Operations

- Element-wise add
- Element-wise plus
- Element-wise division
- Hadamard product
- Matrix multiplication
- Batched matrix multiplication
- More linear algebra operations...
- Collect, Scatter, Reduce...

Libraries

- Numpy
- Blas
- cuBlas
- cuSparse
- MKL
- TensorFlow
- PyTorch
- PaddlePaddle
- MXNet
- ...

Lazy Evaluation and Code Generation

$c = a + b$

$d = c * 2$

for $i = 1$ to n do

$c[i] = a[i] + b[i]$

for $i = 1$ to n do

$d[i] = c[i] * 2$

for $i = 1$ to n do

$d[i] = (a[i] + b[i]) * 2$

Optimizations

- Graph minimization and canonicalization
 - Constant Folding
 - Common subexpression elimination
 - Remove unnecessary operations
- Algebraic simplification and reassociation
- Copy propagation

Graph Optimizer

- Graph minimization and canonicalization
 - Constant Folding
 - Common subexpression elimination
 - Remove unnecessary operations
- Algebraic simplification and reassociation
- Copy propagation

Meta Optimizer

```
i = 0
while i < config.meta_optimizer_iterations (default=2):
    Pruning ()           # Remove nodes not in fanin of outputs, unused functions
    Function ()         # Function specialization & inlining, symbolic gradient inlining
    DebugStripper () *  # Remove assert, print, check_numerics
    ConstFold ()       # Constant folding and materialization
    Shape ()           # Symbolic shape arithmetic
    Remapper ()        # Op fusion
    Arithmetic ()     # Node deduping (CSE) & arithmetic simplification
    if i==0: Layout () # Layout optimization for GPU
    if i==0: Memory () # Swap-out/Swap-in, Recompute*, split large nodes
    Loop ()           # Loop Invariant Node Motion*, Stack Push & Dead Node Elimination
    Dependency ()     # Prune/optimize control edges, NoOp/Identity node pruning
    Custom ()         # Run registered custom optimizers (e.g. TensorRT)
    i += 1
```

Constant Folding Optimizer

```
do:  
  InferShapesStatically() # Fixed-point iteration with symbolic shapes  
  graph_changed = MaterializeConstants() # grad broadcast, reduction dims  
  q = NodesWithKnownInputs()  
  while not q.empty():  
    node = q.pop()  
    graph_changed |= FoldGraph(node, &q) # Evaluate node on host  
  graph_changed |= SimplifyGraph()  
while graph_changed
```

Constant Folding Optimizer: SimplifyGraph()

- Removes trivial ops, e.g. identity Reshape, Transpose of 1-d tensors, $\text{Slice}(x) = x$, etc.
- Rewrites that enable further constant folding
- Arithmetic rewrites that rely on known shapes or inputs, e.g.
 - Constant push-down:
 - $\text{Add}(c1, \text{Add}(x, c2)) \Rightarrow \text{Add}(x, c1 + c2)$
 - $\text{ConvND}(c1 * x, c2) \Rightarrow \text{ConvND}(x, c1 * c2)$
 - Partial constfold:
 - $\text{AddN}(c1, x, c2, y) \Rightarrow \text{AddN}(c1 + c2, x, y)$,
 - $\text{Concat}([x, c1, c2, y]) = \text{Concat}([x, \text{Concat}([c1, c2]), y])$
 - Operations with neutral & absorbing elements:
 - $x * \text{Ones}(s) \Rightarrow \text{Identity}(x)$, if $\text{shape}(x) == \text{output_shape}$
 - $x * \text{Ones}(s) \Rightarrow \text{BroadcastTo}(x, \text{Shape}(s))$, if $\text{shape}(s) == \text{output_shape}$
 - Same for $x + \text{Zeros}(s)$, $x / \text{Ones}(s)$, $x * \text{Zeros}(s)$ etc.
 - $\text{Zeros}(s) - y \Rightarrow \text{Neg}(y)$, if $\text{shape}(y) == \text{output_shape}$
 - $\text{Ones}(s) / y \Rightarrow \text{Recip}(y)$ if $\text{shape}(y) == \text{output_shape}$

Arithmetic Optimizer

```
DedupComputations():
```

```
do:
```

```
    stop = true
```

```
    UniqueNodes reps
```

```
    for node in graph.nodes():
```

```
        rep = reps.FindOrInsert(node, IsCommutative(node))
```

```
        if rep == node or !SafeToDedup(node, rep):
```

```
            continue
```

```
        for fanout in node.fanout():
```

```
            ReplaceInputs(fanout, node, rep)
```

```
        stop = false
```

```
while !stop
```

Arithmetic Optimizer

- Arithmetic simplifications
 - Flattening: $a+b+c+d \Rightarrow \text{AddN}(a, b, c, d)$
 - Hoisting: $\text{AddN}(x * a, b * x, x * c) \Rightarrow x * \text{AddN}(a+b+c)$
 - Simplification to reduce number of nodes:
 - Numeric: $x+x+x \Rightarrow 3*x$
 - Logic: $!(x > y) \Rightarrow x \leq y$
- Broadcast minimization
 - Example: $(\text{matrix1} + \text{scalar1}) + (\text{matrix2} + \text{scalar2}) \Rightarrow (\text{matrix1} + \text{matrix2}) + (\text{scalar1} + \text{scalar2})$
- Better use of intrinsics
 - $\text{Matmul}(\text{Transpose}(x), y) \Rightarrow \text{Matmul}(x, y, \text{transpose_x}=\text{True})$
- Remove redundant ops or op pairs
 - $\text{Transpose}(\text{Transpose}(x, \text{perm}), \text{inverse_perm})$
 - $\text{BitCast}(\text{BitCast}(x, \text{dtype1}), \text{dtype2}) \Rightarrow \text{BitCast}(x, \text{dtype2})$
 - Pairs of elementwise involutions $f(f(x)) \Rightarrow x$ (Neg, Conj, Reciprocal, LogicalNot)
 - Repeated Idempotent ops $f(f(x)) \Rightarrow f(x)$ (DeepCopy, Identity, CheckNumerics...)
- Hoist chains of unary ops at Concat/Split/SplitV
 - $\text{Concat}([\text{Exp}(\text{Cos}(x)), \text{Exp}(\text{Cos}(y)), \text{Exp}(\text{Cos}(z))]) \Rightarrow \text{Exp}(\text{Cos}(\text{Concat}([x, y, z])))$
 - $[\text{Exp}(\text{Cos}(y)) \text{ for } y \text{ in Split}(x)] \Rightarrow \text{Split}(\text{Exp}(\text{Cos}(x)), \text{num_splits})$

Tensor

- Dense
 - Column major
 - Row major
 - Stride
- Sparse
 - Compressed representation
 - Set intersection

Differentiation

- Numerical differentiation
- Symbolic differentiation
 - Chain rules
- Forward mode auto differentiation
- Reverse mode auto differentiation